

**ANALYZING PRIMARY STUDENT DATA USING DATA
MINING**

CHONG SZE WEI

**UNIVERSITI UTARA MALAYSIA
2009**

ANALYZING PRIMARY STUDENT DATA USING DATA MINING

A dissertation submitted to the Division of Applied Sciences, College of Arts and Sciences in partial fulfillment of the requirements for the degree Master of Science (Information Technology), Universiti Utara Malaysia.

By

Chong Sze Wei

© Chong Sze Wei, 2009. All rights reserved.

PERMISSION TO USE

In presenting this dissertation in partial fulfillment of the requirements for a postgraduate degree from Universiti Utara Malaysia, I agree that the University Library may make it freely available for inspection. I further agree that permission for copying of this dissertation in any manner, in whole or in part, for scholarly purpose may be granted by my supervisor(s) or, in their absence by the Dean of College of Arts and Sciences. It is understood that any copying or publication or use of this dissertation or parts thereof for financial gain shall not be allowed without my written permission. It is also understood that due recognition shall be given to me and to Universiti Utara Malaysia for any scholarly use which may be made of any material from my dissertation.

Requests for permission to copy or to make other use of materials in this dissertation, in whole or in part, should be addressed to

Division of Applied Sciences

College of Arts and Sciences

Universiti Utara Malaysia

06010 UUM Sintok

Kedah Darul Aman.

ABSTRACT (ENGLISH)

Nowadays, academic achievement has become the most important evidence for establishing the value of Malaysia's education boundary. In this study, the primary students' examination data is collected on the previous examination mark yet sake to be analyzed for their future study plan. The selection of using data mining approaches was based on the capability of data mining as a grateful tool for academic analysis purposes. Focused on educational boundary, data mining approaches can be used for the process of uncovering hidden information and patterns that can help school community forecast the students' academic achievement. Therefore, the other relevant data such as student performance information and family income also engaged in this study. The overall relevant raw datasets is used for preprocessed and analyzed using statistical method. In addition, the result from the statistical manner analysis point out the considerable contribution of these attributes to the academic achievement plan.

ABSTRACT (MALAY)

Saban hari ini, pencapaian dalam bidang akademik sudah menjadi suatu tanda yang kukuh untuk menunjukkan nilai dalam batasan pendidikan Malaysia. Merujuk kepada penyelidikan ini, data peperiksaan yang lepas bagi pelajar pada peringkat sekolah rendah akan dikumpulkan untuk dianalisis sebagai rujukan pelan pelajaran mereka pada masa hadapan. Pemilihan menggunakan teknik data mining adalah ikutan kepada kebolehan Data Mining sebagai satu alatan yang canggih untuk kegunaan bagi analisis tentang bidang akademik. Dengan menumpu kepada batasan pendidikan, teknik data mining boleh digunakan untuk menunjukkan informasi yang tersembunyi dan untuk menunjukkan pola-pola atau corak yang boleh membantu komuniti sekolah menerokai pencapaian pelajar dalam bidang akademik. Justeru, data-data yang berkaitan seperti data kemahiran pelajar dan data latar belakang keluarga juga akan tergolong dalam penyelidikan ini. Secara keseluruhan, data mentah akan digunakan untuk pemprosesan dan kaedah Statistik akan digunakan untuk analisis. Tambahan pula, keputusan yang telah diperolehi akan menunjukkan sumbangan yang boleh ditimbang tentang hal-hal yang berkaitan dengan pelan kecemerlangan dalam bidang akademik.

ACKNOWLEDGEMENT

First and foremost, I would like to express my deepest gratitude to my supervisor, Associate Professor Fadzilah Siraj, who willing to supervise and lead me patiently with plot a routing and charitably sahiring her abundant source of knowledge in this dissertation. Indeed, I really grateful for my supervisor because without her assistance numerous of beneficial comments, this study would have never been possible.

Therefore, I wish to send my warmest gratitude to my beloved parents, Mr. Chong Chan Mui and Mrs. Yoo Bee Eng for all their love and support by furnishing me a stunning position in their heart. In addition, thankful to my siblings which includes Chong Sze Hui, Chong Sien Joo, and Chong Sze Yin respectively for their help and encouragement.

Last but not least, a million thanks to all those who had lent a helping hand in permitting me to become visible of this project.

Chong Sze wei

Division of Applied Sciences

College of Arts and Sciences

TABLE OF CONTENTS

Permission to Use.....	i
Abstract (English).....	ii
Abstract (Malay).....	iii
Acknowledgement.....	iv
Table of Contents.....	v
List of Figures.....	viii
List of Tables.....	xii
List of Abbreviations.....	xiii

CHAPTER 1: INTRODUCTION

1.1 Background of Study.....	1
1.2 Problem Statements.....	3
1.3 Research Objectives.....	4
1.4 Research Scope.....	4
1.5 Significance of the Study.....	5
1.6 Structure of Report.....	5

CHAPTER 2: LITERATURE REVIEW

2.1 Data Mining.....	6
2.2 Descriptive of Data Mining.....	7
2.3 Predictive of Data Mining.....	9

2.4 Data Mining in Education.....	10
2.5 Benefits of Data Mining.....	15

CHAPTER 3: RESEARCH METHODOLOGY

3.1 Integrated Data.....	18
3.1.1 Instrument.....	19
3.1.2 Respondents Data.....	20
3.1.3 Target.....	20
3.2 Cross Industry Standard Process.....	21
3.3 Information System Development Research Process.....	25

CHAPTER 4: FINDINGS AND RESULTS

4.1 Prototype of BICC	29
4.2 Descriptive Analysis	31
4.2.1 Year 4.....	31
4.2.2 Year 5.....	52
4.2.3 Year 6.....	72

CHAPTER 5: CONCLUSION AND RECOMMENDATIONS

5.1 Conclusion.....	94
5.2 Recommendations.....	95

REFERENCES

APPENDICES

Appendix A: Raw Dataset

Appendix B: Descriptive Analysis

Appendix C: User Manual

TABLE OF FIGURES

Figure	Title	Page
3.1	Process flow of the study.....	17
3.2	Snapshots of the upload file screen.....	19
3.3	Phase of CRISP-DM Reference Model.....	21
3.4	Sample of raw data set before converted using SPSS 12.0.....	24
3.5	Process of data set conversion using SPSS 12.0.....	24
3.6	Data set after converted in SPSS 12.0.....	25
3.7	Information System Development Research Process Model.....	26
3.8	Prototype System Structure Diagram.....	27
4.1	Snapshot of the login page	29
4.2	Snapshot of the upload file page	30
4.3	Snapshot of the analysis outcome page	30
4.4	Gender distribution.....	32
4.5	Gender versus BM1.....	32
4.6	Gender versus BM2.....	33
4.7	Gender versus ENG.....	34
4.8	Gender versus MATH.....	34
4.9	Gender versus Science.....	35
4.10	Attendance distribution.....	36
4.11	Attendance versus BM1.....	37
4.12	Attendance versus BM2.....	38

4.13	Attendance versus ENG.....	38
4.14	Attendance versus MATH.....	39
4.15	Attendance versus Science.....	40
4.16	Co-curricular activities distribution.....	41
4.17	Co-curricular activities versus BM1.....	41
4.18	Co-curricular activities versus BM2.....	42
4.19	Co-curricular activities versus ENG.....	43
4.20	Co-curricular activities versus MATH.....	44
4.21	Co-curricular activities versus Science.....	45
4.22	Family income distribution.....	46
4.23	Family income versus BM1.....	47
4.24	Family income versus BM2.....	48
4.25	Family income versus ENG.....	49
4.26	Family income versus MATH.....	50
4.27	Family income versus Science.....	51
4.28	Gender distribution.....	53
4.29	Gender versus BM1.....	53
4.30	Gender versus BM2.....	54
4.31	Gender versus ENG.....	55
4.32	Gender versus MATH.....	56
4.33	Gender versus Science.....	56
4.34	Attendance distribution.....	57
4.35	Attendance versus BM1.....	58

4.36	Attendance versus BM2.....	58
4.37	Attendance versus ENG.....	59
4.38	Attendance versus MATH.....	60
4.39	Attendance versus Science.....	61
4.40	Co-curricular activities distribution.....	62
4.41	Co-curricular activities versus BM1.....	62
4.42	Co-curricular activities versus BM2.....	63
4.43	Co-curricular activities versus ENG.....	64
4.44	Co-curricular activities versus MATH.....	65
4.45	Co-curricular activities versus Science.....	66
4.46	Family income distribution.....	67
4.47	Family income versus BM1.....	67
4.48	Family income versus BM2.....	68
4.49	Family income versus ENG.....	69
4.50	Family income versus MATH.....	70
4.51	Family income versus Science.....	71
4.52	Gender distribution.....	73
4.53	Gender versus BM1.....	73
4.54	Gender versus BM2.....	74
4.55	Gender versus ENG.....	75
4.56	Gender versus MATH.....	76
4.57	Gender versus Science.....	76
4.58	Attendance distribution.....	77

4.59	Attendance versus BM1.....	78
4.60	Attendance versus BM2.....	79
4.61	Attendance versus ENG.....	80
4.62	Attendance versus MATH.....	81
4.63	Attendance versus Science.....	82
4.64	Co-curricular activities distribution.....	83
4.65	Co-curricular activities versus BM1.....	83
4.66	Co-curricular activities versus BM2.....	84
4.67	Co-curricular activities versus ENG.....	85
4.68	Co-curricular activities versus MATH.....	86
4.69	Co-curricular activities versus Science.....	87
4.70	Family income distribution.....	88
4.71	Family income versus BM1.....	89
4.72	Family income versus BM2.....	90
4.73	Family income versus ENG.....	91
4.74	Family income versus MATH.....	92
4.75	Family income versus Science.....	93

TABLE OF TABLES

Table	Title	Page
3.1	Examples of records attributes.....	18
3.2	Grade categorized.....	20
3.3	Respondent from various class.....	22
3.4	The selected attributes before converted to numeric.....	23
3.5	The selected attributes after converted to numeric.....	23
4.1	Frequency data.....	31
4.2	Frequency data.....	52
4.3	Frequency data.....	72

TABLE OF ABBREVIATIONS

ASP	Active Server Page
BICC	Business Intelligence Competency Centre
BM1	Bahasa Melayu 1
BM2	Bahasa Melayu 2
DM	Data Mining
DMS	Data Mining System
ENG	English
GEP	Gifted Education Programme
ITS	Intelligent Tutoring System
KDD	Knowledge Discovery in Database
MATH	Mathematic
MOE	Ministry of Education
PMR	Penilaian Menengah Rendah
RM	Ringgit Malaysia
SMM	Sistem Maklumat Murid
SPM	Sijil Pelajaran Malaysia
SRKBG	Sekolah Rendah Kebangsaan Griik
STPM	Sijil Tinggi Pelajaran Malaysia
UPSR	Ujian Penilaian Sekolah Rendah
USA	United State

CHAPTER 1

INTRODUCTION

This chapter discusses the background of the study that consists of several sub-parts about subject of the study. These include overview on data mining techniques used for education area, a short description on the problem statements, research objectives, research scope and significance of this study. Lastly, this chapter presents the study organization by describing the structure of this report.

1.1 Background of Study

Data mining (DM) can be defined as the processes of extracting interesting information such as non-trivial, implicit, potentially valuable or previously unknown information from a huge amount of data storage area such as data warehouse, relational database and so on (Chen *et al.* 1996). It is a very valuable way of analyzing a huge amount of data, especially when humans are not capable to analyze such datasets manually. Also DM is well known as one of the core processes of Knowledge Discovery in Database (KDD).

DM also is known as one of the key features of many homeland safety initiatives. Regularly used as a means for detecting fraud, product retailing and others, data

The contents of
the thesis is for
internal user
only

REFERENCES

- Anjewierden, A., Koll'offel, B. & Hulshof, C. (2007). Using data mining methods for automated chat analysis to understand and support inquiry learning processes. *Proceeding of Towards Educational Data Mining*. Enschede, The Netherlands. pp. 27-36.
- Arnold, A., Beck, J. & Scheines, R. (2004). An Inductive Approach. *Proceeding of Feature Discovery in the Context of Educational Data Mining*. Pittsburgh, PA 15213, USA.
- Beikzadeh, M. R. & Delavari, N. (2004). A new analysis model for data mining processes in higher educational systems. *Proceeding of 5th international conference ITHET*. MMU, Cyberjaya, Malaysia. pp. 5-8.
- Beal, C. R. & Cohen, P. R. (2006). Temporal Data Mining for Educational Applications.
- Bravo, J., Vialardi, C. & Ortigosa, A. (2007). A problem-oriented method for supporting AEH authors through data mining. *Proceeding of International Workshop on Applying Data Mining in e-Learning (ADML'07)*. Madrid, Spain. pp. 53-62.
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C. & Wirth, R. (1999). CRISP-DM 1.0. *Proceeding of Step-by-step Data Mining Guide*. pp. 9-10.
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C. & Wirth, R. (1999). CRISP-DM 1.0. *Proceeding of Step-by-step Data Mining Guide*. pp. 9-10.
- Chen, M.-S., Jan, J. & Yu, P.S. (1996). Data mining: An overview from a database perspective. *Proceeding of IEEE Transaction on Knowledge and Data Engineering*, 8. pp. 866-883.

- Garcia, E., Romero, C., Ventura, S. & Calder, T. (2007). Drawbacks and solutions of applying association rule mining in learning management systems. *Proceeding of International Workshop on Applying Data Mining in e-Learning (ADML'07)*. pp. 13-22.
- Gargano, M. L. & Raggad, B. G. (1999). Data mining-a powerful information creating tool. *Proceeding of OCLC Systems & Services*. 15(2), pp. 81-90.
- Hamalainen, W., Laine, T. H. & Suitinen, E. (2004). *Data Mining in Personalizing Distance Education Courses*. University of Joensuu, Finland.
- Lloyd, N. M., Heffernan, N. T. & Ruiz, C. (2007). Predicting student engagement in intelligent tutoring systems using teacher expert knowledge. *Proceeding of Educational Data Mining Workshop*. Marina del Rey, CA. USA. pp. 40-49.
- Luan, J. (2001). Data mining as driven by knowledge management in higher education: Persistence clustering and prediction. *Proceeding of Keynote for SPSS Public Cpnference, UCSF*. pp. 1-16.
- Ma, Y., Liu, B., Wong, C. K., Yu, P. S. & Lee, S. M. (2000). Targeting the Right Students Using Data Mining.
- Merceron, A. & Yacef, K. (2007). Revisiting interestingness of strong symmetric association rules in education data. *Proceeding of International Workshop on Applying Data Mining in e-Learning (ADML'07)*. pp. 3-12.
- Minaei-Bidgoli, B., Kortemeyer, G. & Punch, W. F. (2003). Optimizing classification ensembles via a genetic algorithm for a web-based educational system. East Lansing, USA. pp.1-9.

- Minaei-Bidgoli, B., Tan, P. N. & Punch, W. F. (2004). Mining interesting contrast rules for a web-based educational system. East Lansing, USA. pp. 1-8.
- Perera, D., Kay, J., Yacef, K. & Koprinska I. (2007). Mining learners' traces from an online collaboration tool. *Proceeding of Educational Data Mining Workshop*. University of Sydney, Australia. pp. 60-69.
- Ravi, S., Kim, J. & Shaw, E. (2007). Mining on-line discussions: Assessing technical quality for student scaffolding and classifying messages for participation profiling. *Proceeding of Educational Data Mining Workshop*. Marina del rey, CA. USA. pp. 70-79.
- Raghavan, V. & Hafez, A. (2000). *Dynamic Data Mining*. University of Louisiana, USA. pp. 1-10.
- Reza, B. & Naeimeh, D. (2004). *A New Analysis Model for Data Mining Processes in Higher Educational Systems*. MMU, Cyberjaya. Malaysia. pp. 5-8.
- Seifert, J. W. (2004). Data mining: An overview. *Proceeding of CRS Report for Congress*. pp. 1-16.
- Tanimoto, S. T. (2007). Improving the prospects for educational data mining. Pp. 1-6.
- Tsantis, L. & Castellani, J. (2001). Enhancing learning environments through solution-based knowledge discovery tools: Forecasting for self-perpetuating systemic reform. pp. 1-35.
- Yudelson, M. V., Medvedeva, O., Legowski, E., Castine, M., Jukic, D. & Crowley, R. S. (2006). Mining student learning data to develop high level pedagogic strategy in a medical ITS. University of Pittsburgh, pp. 1-8.